



ELSEVIER

Theoretical Computer Science 284 (2002) 181–196

Theoretical
Computer Science

www.elsevier.com/locate/tcs

On the rate of convergence of error estimates for the partitioning classification rule[☆]

Márta Pintér¹

*Department of Computer Science and Information Theory, Technical University of Budapest,
1521 Stoczek u. 2, Budapest, Hungary*

Abstract

The error probability of the partitioning classification rule is shown to converge to the Bayes error faster than $1/\sqrt{n}$ under certain conditions. The resubstitution and the deleted error estimates for the partitioning classification rule from a sample $(X_1, Y_1), \dots, (X_n, Y_n)$ are studied. The random part of the resubstitution estimate is shown to be small for arbitrary partition and for any distribution of (X, Y) . If we assume that X has a density f and the partitions consist of rectangles, then the difference between the expected value of the estimate and the Bayes error restricted to the partition is less than a constant times $1/\sqrt{n}$. The main result of the paper is that, under the same conditions, for both estimates the difference between the estimate and the real error probability of the classification rule is asymptotically normal with 0 mean and variance $L^*/2$, where L^* is the Bayes error. © 2002 Elsevier Science B.V. All rights reserved.

Keywords: Error estimation; Asymptotic normality; Rate of convergence; Partitioning classification rule

1. Introduction

Let X be the d -dimensional feature vector with distribution μ and let Y be the binary valued label. Denote the posterior probabilities by

$$P_i(x) = \mathbf{P}\{Y = i|X = x\}, \quad i = 0, 1.$$

In pattern recognition the value of the label Y is to be predicted upon observing the feature vector X . The prediction rule or classifier g is a function $\mathcal{R}^d \rightarrow \{0, 1\}$ whose

[☆] Parts of the paper were presented at EUROCOLT'99.

¹ The research of the author was supported by the Computer and Automation Institute of the Hungarian Academy of Sciences (MTA SZTAKI).

E-mail address: marti@szit.bme.hu (M. Pintér).

performance is measured by the probability of error

$$L(g) = \mathbf{P}\{g(X) \neq Y\}.$$

The Bayes classifier

$$g^*(x) = \begin{cases} 0 & \text{if } P_1(x) \leq 1/2, \\ 1 & \text{otherwise.} \end{cases}$$

is well-known to have minimal probability of error among all possible classifiers. Its error probability $L(g^*)$ is called the Bayes error and is denoted by L^* .

$$L^* = \mathbf{E}\{\min(P_0(X), P_1(X))\}.$$

Assume that n independent copies of (X, Y) form the available data sequence:

$$D_n = ((X_1, Y_1), \dots, (X_n, Y_n)).$$

These data may be used to design the classification rule $g_n(x)$ whose probability of error is the random variable

$$L_n = L(g_n) = \mathbf{P}\{g_n(X) \neq Y | D_n\}.$$

Many important classification rules partition \mathcal{R}^d into disjoint cells and classify in each cell according to the majority vote among the labels of the sample points X_i falling in the same cell. Let $\mathcal{P}_n = \{A_{nj}, j = 1, 2, \dots\}$ be a partition of \mathcal{R}^d and let $A_n(x)$ denote the cell in the partition that includes x . Then the partitioning classification rule:

$$g_n(x) = \begin{cases} 0 & \text{if } \sum_{i=1}^n I_{\{Y_i=1\}} I_{\{X_i \in A_n(x)\}} \leq \sum_{i=1}^n I_{\{Y_i=0\}} I_{\{X_i \in A_n(x)\}}, \\ 1 & \text{otherwise.} \end{cases}$$

For cubic partitions, consistency was established under some additional conditions by Glick [11]. Universal consistency follows from the result of Gordon and Olshen [12, 13]. The usual conditions needed for the consistency are that the diameter of the cells tends to zero and the number of cells in a bounded region is small compared to the sample size. Strong universal consistency of the partitioning classification rule was proved by Devroye and Györfi [8]. We will show in Section 2 that, although universal rate of convergence does not exist, the error probability of the partitioning classification rule converges to the Bayes error faster than $1/\sqrt{n}$ under certain conditions. In Section 3 we study two error estimates, the resubstitution and the deleted estimate (or cross validation). Asymptotic normality of the resubstitution estimate was earlier studied in Györfi and Horváth [14]: there it is shown that the difference between the resubstitution estimate and the Bayes error restricted to the partition is asymptotically normal. We prove that the difference between the estimate and the true error probability is asymptotically normal and that the same holds also for the deleted estimate.

We will need in all results the assumption that the labels in a cell are not uniformly chosen; more precisely we will assume that there is a constant c such that

$$|\mu(A_{nj}) - 2v(A_{nj})| > c\mu(A_{nj}) \quad (1)$$

for all n and j , where $v(A) = \mathbf{E}\{YI_{\{X \in A\}}\}$.

The necessity of such a condition is demonstrated in Proposition 1, where we show that if the label is uniform in every cell then the difference between the resubstitution estimate and the true error probability is too big and asymptotic normality does not hold.

Partitioning estimates, which were introduced by Tukey [16] for regression estimation, are present also in density estimation; their behaviour has been studied in several papers and they are known to be consistent. Asymptotic normality of the L_1 - and L_2 -errors in partitioning density estimation (histogram) was proved by Beirlant et al. [3] and Berinet et al. [5]. Beirlant and Mason [4] developed a general method for deriving asymptotic normality of the L_p -errors of empirical functionals which make use of the neighboring data at any point x of interest. Beirlant and Györfi [2] continued this line of research and established the asymptotic normality of L_2 -error in partitioning regression estimation. To obtain asymptotic normality all these authors used the idea of Poissonization in conjunction with Bartlett's [1] method of partial inversion for obtaining characteristic functions of conditional distributions.

In the paper, we assume that X takes its value in a bounded region $\mathcal{X} \subset \mathcal{R}^d$.

2. Rate of convergence

It is well-known that no universal rate of convergence can be guaranteed for any classifiers if nothing is assumed about the distribution. There is no sequence of numbers a_n converging to zero such that there is a classification rule with error probability below $L^* + a_n$ for all distributions [7] (see also [10, Section 7.2]). In this section we show that the error probability of the partitioning classification rule converges to the Bayes error faster than $1/\sqrt{n}$ under certain conditions.

Theorem 1. *Let R_n^* denote the Bayes error restricted to the partition \mathcal{P}_n , i.e.*

$$R_n^* = \sum_{j=1}^{m_n} \min\{v(A_{nj}), \mu(A_{nj}) - v(A_{nj})\}.$$

Then

- (i) $L_n \geq R_n^*$.
- (ii) *Assume that the partitions are finite, let m_n denote the number of cells in the partition. If (1) is satisfied, then*

$$\mathbf{E}(L_n - R_n^*) \leq \frac{12}{ce} \frac{m_n}{n}.$$

Proof. (i) Clearly $L_n \geq R_n^*$ because R_n^* is the Bayes error corresponding to the pair of random variables $(T_n(X), Y)$, where $T_n(x) = j$ if $x \in A_{nj}$ and since the partitioning classification rule $g_n(x)$ can be written as a function of $T_n(x)$, its error probability L_n cannot be smaller than R_n^* .

(ii) The decision of the partitioning rule in a cell is different from the Bayes decision exactly if either there are more ones in the cell than zeros but $v(A_{nj})$, the probability that a label in the cell is 1 is less than $\mu(A_{nj}) - v(A_{nj})$, the probability that the label is 0 or vice versa. And in that case the difference between L_n and the Bayes error is the difference between the two probabilities: $|\mu(A_{nj}) - v(A_{nj}) - v(A_{nj})| = |\mu(A_{nj}) - 2v(A_{nj})|$.

$$\mathbf{E}(L_n - R_n^*) = \mathbf{E} \left(\sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| I_{\{\text{sign}(\mu_n(A_{nj}) - 2v_n(A_{nj})) \neq \text{sign}(\mu(A_{nj}) - 2v(A_{nj}))\}} \right),$$

where

$$v_n(A) = \frac{1}{n} \sum_{i=1}^n I_{\{X_i \in A\}} Y_i \quad \text{and} \quad \mu_n(A) = \frac{1}{n} \sum_{i=1}^n I_{\{X_i \in A\}}.$$

Thus

$$\begin{aligned} & \mathbf{E}(L_n - R_n^*) \\ &= \mathbf{E} \left(\sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| I_{\{\text{sign}(\mu_n(A_{nj}) - 2v_n(A_{nj})) \neq \text{sign}(\mu(A_{nj}) - 2v(A_{nj}))\}} \right) \\ &= \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| \mathbf{P}(\text{sign}(\mu_n(A_{nj}) - 2v_n(A_{nj})) \neq \text{sign}(\mu(A_{nj}) - 2v(A_{nj}))) \\ &= \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| [\mathbf{P}(\mu_n(A_{nj}) - 2v_n(A_{nj}) > 0) I_{\{\mu(A_{nj}) - 2v(A_{nj}) \leq 0\}} \\ &\quad + \mathbf{P}(\mu_n(A_{nj}) - 2v_n(A_{nj}) \leq 0) I_{\{\mu(A_{nj}) - 2v(A_{nj}) > 0\}}] \\ &= \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| \\ &\quad \times [\mathbf{P}(\mu_n(A_{nj}) - 2v_n(A_{nj}) - (\mu(A_{nj}) - 2v(A_{nj})) > -(\mu(A_{nj}) - 2v(A_{nj}))) \\ &\quad \times I_{\{\mu(A_{nj}) - 2v(A_{nj}) \leq 0\}} \\ &\quad + \mathbf{P}(\mu_n(A_{nj}) - 2v_n(A_{nj}) - (\mu(A_{nj}) - 2v(A_{nj})) \leq -(\mu(A_{nj}) - 2v(A_{nj}))) \\ &\quad \times I_{\{\mu(A_{nj}) - 2v(A_{nj}) > 0\}}] \\ &\leq \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| \\ &\quad \times [\mathbf{P}(\mu_n(A_{nj}) - 2v_n(A_{nj}) - (\mu(A_{nj}) - 2v(A_{nj})) > |\mu(A_{nj}) - 2v(A_{nj})|) \\ &\quad + \mathbf{P}(\mu_n(A_{nj}) - 2v_n(A_{nj}) - (\mu(A_{nj}) - 2v(A_{nj})) \leq -|\mu(A_{nj}) - 2v(A_{nj})|)] \\ &\leq 2 \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| \\ &\quad \times \mathbf{P}(|\mu_n(A_{nj}) - 2v_n(A_{nj}) - (\mu(A_{nj}) - 2v(A_{nj}))| > |\mu(A_{nj}) - 2v(A_{nj})|). \quad (2) \end{aligned}$$

We can apply Lemma A.1 of the appendix with $\varepsilon = |\mu(A) - 2v(A)|$ and $Z_i = I_{\{X_i \in A\}}(1 - 2Y_i) - (\mu(A) - 2v(A))$ because they are independent identically distributed zero mean random variables and $|Z_i| \leq 2$:

$$\begin{aligned} \mathbf{P}(|\mu_n(A) - 2v_n(A) - (\mu(A) - 2v(A))| > |\mu(A) - 2v(A)|) \\ \leq 2 \exp \left(\frac{-n|\mu(A) - 2v(A)|^2}{2(\sigma_A^2 + 2|\mu(A) - 2v(A)|)} \right), \end{aligned}$$

where $\sigma_A^2 = \mathbf{E}(Z_i^2) = \mu(A) - (\mu(A) - 2v(A))^2$.

From this and (2)

$$\begin{aligned} \mathbf{E}(L_n - R_n^*) &\leq 2 \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| 2e^{\frac{-n|\mu(A_{nj}) - 2v(A_{nj})|^2}{2(\mu(A_{nj}) - (\mu(A_{nj}) - 2v(A_{nj}))^2 + 2|\mu(A_{nj}) - 2v(A_{nj})|)}} \\ &\leq 2 \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| 2e^{\frac{-n|\mu(A_{nj}) - 2v(A_{nj})|^2}{2(3\mu(A_{nj}))}} \\ &\leq 2 \sum_{j=1}^{m_n} |\mu(A_{nj}) - 2v(A_{nj})| 2e^{-\frac{c}{6}n|\mu(A_{nj}) - 2v(A_{nj})|} \\ &\quad \text{(because of condition (1))} \\ &= 2 \frac{6}{cn} \sum_{j=1}^{m_n} \frac{c}{6} n |\mu(A_{nj}) - 2v(A_{nj})| 2e^{-\frac{c}{6}n|\mu(A_{nj}) - 2v(A_{nj})|} \\ &\leq 4 \frac{6}{cn} m_n \max_z z e^{-z} \\ &= \frac{12}{ce} \frac{m_n}{n}. \quad \square \end{aligned}$$

Corollary 1. *If in addition to condition (1)*

$$\lim_{n \rightarrow \infty} \frac{m_n^2}{n} = 0, \tag{3}$$

then

$$\sqrt{n}(L_n - R_n^*) \rightarrow 0$$

in probability as $n \rightarrow \infty$.

Proof. It is well-known that if, for a nonnegative random variable Z_n , $\mathbf{E}Z_n \rightarrow 0$, then $Z_n \rightarrow 0$ in probability. Thus we have to prove that

$$\sqrt{n}\mathbf{E}(L_n - R_n^*) \rightarrow 0$$

as $n \rightarrow \infty$. From Theorem 1

$$\sqrt{n}\mathbf{E}(L_n - R_n^*) \leq \frac{12e^{-1}}{c} \frac{m_n}{\sqrt{n}} \rightarrow 0$$

because of condition (3). \square

3. Error estimation

To derive two important corollaries of Theorem 1, we introduce two error estimates: the resubstitution error estimate and the deleted estimate (or cross-validation).

Estimating the error probability of a classification rule g_n is of great importance. The designer always wants to know what performance can be expected from a classifier. Since the distribution of the data is unknown, it is important to find and analyze error estimation methods that work well independently from the distribution of (X, Y) .

One natural way of error estimation is to split the data into two parts: a training sample and a testing sample, and use the first for designing a classifier and the latter for the estimation. But when data are expensive this constitutes a waste. There are well-known methods that use the same sample for training and testing: the resubstitution estimate, the deleted estimate, etc.

3.1. Resubstitution estimate

The resubstitution estimate \hat{L}_n counts the number of errors committed on the training sequence $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ by the classification rule, i.e. for a classifier g_n it is defined as

$$\hat{L}_n = \frac{1}{n} \sum_{i=1}^n I_{\{g_n(X_i) \neq Y_i\}},$$

which for the partitioning classification rule can be written in the following form:

$$\hat{L}_n = \sum_{j=1}^{m_n} \min\{v_n(A_{nj}), \mu_n(A_{nj}) - v_n(A_{nj})\}, \quad (4)$$

where

$$v_n(A) = \frac{1}{n} \sum_{i=1}^n I_{\{X_i \in A\}} Y_i \quad \text{and} \quad \mu_n(A) = \frac{1}{n} \sum_{i=1}^n I_{\{X_i \in A\}}.$$

\hat{L}_n is an estimate of the Bayes error R_n^* restricted to the partition \mathcal{P}_n .

Concerning the resubstitution error estimate for partitioning rule the following inequalities are known (see [10, Section 23.2]): for an arbitrary partition \mathcal{P}_n

$$\text{Var}(\hat{L}_n) \leq \frac{1}{n}$$

and

$$\mathbf{E}\hat{L}_n \leq R_n^*.$$

For a finite partition of size m_n

$$R_n^* - \mathbf{E}\hat{L}_n \leq \sqrt{\frac{2m_n}{n}}. \quad (5)$$

Combining this with the result of Theorem 1, we get

$$\mathbf{E}(L_n - \hat{L}_n) \leq \sqrt{\frac{2m_n}{n}} + \frac{12}{ce} \frac{m_n}{n}.$$

The resubstitution estimate for the partitioning rule is asymptotically normal under certain conditions:

Theorem 2. Consider the partitions where A_{nj} are d -dimensional rectangles. Let $a_n^i(x)$, $i = 1, 2, \dots, d$, denote the sidelengths of $A_n(x)$. Assume that for all x there exists a $K(x)$ so that $a_n^i(x)/a_n^j(x) \leq K(x)$ for all $1 \leq i, j \leq d$ and n . Assume that

$$\lim_{n \rightarrow \infty} \sup_{A_{nj}} \text{diam}(A_{nj}) = 0$$

and

$$\lim_{n \rightarrow \infty} \frac{\log n}{n\lambda(A_n(x))} = 0,$$

where λ is the Lebesgue measure. If μ has a density and conditions (1) and (3) are satisfied, then

$$n^{1/2}(L_n - \hat{L}_n)/\sqrt{L^*/2} \xrightarrow{\mathcal{D}} N(0, 1).$$

Thus the resubstitution estimate is a good estimate of the true error probability of the partitioning classification rule and using the result of the theorem a good estimate can be given to the probability that the difference between the resubstitution estimate and the error probability is greater than ε/\sqrt{n} . The result is a consequence of Corollary 1 and Lemma A.2 of the appendix.

Remark. For cubic partitions with size h_n these conditions mean $h_n \rightarrow 0$ and $nh_n^{2d} \rightarrow \infty$ as $n \rightarrow \infty$. Although the second condition is stronger than the condition $nh_n^d \rightarrow \infty$ for the consistency of the partitioning rule, it does allow the usual good choices for h_n , i.e. choices that minimize the rate of convergence, e.g. $h_n = n^{-1/3}$. The condition might be weakened with a longer and more complicated proof.

Under the conditions of Theorem 2 we can remove m_n from the upper bound in (5).

Theorem 3. For any distribution of (X, Y) and for n large enough for the estimate \hat{L}_n of the error probability of a partitioning rule satisfying the conditions of Theorem 2

$$R_n^* - \mathbf{E}\hat{L}_n \leq \frac{1.18}{\sqrt{n}}.$$

Proof. From the result on the asymptotic normality

$$\mathbf{P}\{n^{1/2}|\hat{L}_n - R_n^*| > \varepsilon\} \rightarrow 1 - 2\Phi\left(\frac{\varepsilon}{\sqrt{L^*/2}}\right),$$

where Φ is the standard normal distribution function. Since $R_n^* - \mathbf{E}\hat{L}_n$ is not random

$$\begin{aligned} I_{\{\sqrt{n}|R_n^* - \mathbf{E}\hat{L}_n| > \varepsilon\}} &= \mathbf{P}\{\sqrt{n}|R_n^* - \mathbf{E}\hat{L}_n| > \varepsilon\} \\ &\leq \mathbf{P}\{\sqrt{n}(|\hat{L}_n - \mathbf{E}\hat{L}_n| + |\hat{L}_n - R_n^*|) > \varepsilon\} \\ &\leq \mathbf{P}\{\sqrt{n}|\hat{L}_n - \mathbf{E}\hat{L}_n| > \varepsilon/2\} \\ &\quad + \mathbf{P}\{\sqrt{n}|\hat{L}_n - R_n^*| > \varepsilon/2\}. \end{aligned}$$

For the first term we can apply McDiarmid's inequality (Lemma A.3) with $Z_i = (X_i, Y_i)$ and $A = \mathcal{P}^d \times \{0, 1\}$. It can be easily seen that in the case of the resubstitution estimate for partitioning rule $c_i = 1/n$. Thus for all n and $\varepsilon > 0$,

$$\mathbf{P}\{n^{1/2}|\hat{L}_n - \mathbf{E}\hat{L}_n| > \varepsilon\} \leq 2e^{-2\varepsilon^2}.$$

Therefore

$$\limsup_{n \rightarrow \infty} I_{\{\sqrt{n}|R_n^* - \mathbf{E}\hat{L}_n| > \varepsilon\}} \leq 2e^{-\varepsilon^2/2} + 1 - 2\Phi\left(\frac{\varepsilon}{\sqrt{2L^*}}\right).$$

Obviously there exists such ε_0 for which the right-hand side of the inequality is strictly smaller than 1 and then

$$I_{\{\sqrt{n}|R_n^* - \mathbf{E}\hat{L}_n| > \varepsilon_0\}} = 0$$

for sufficiently large n . Using the trivial upper bound

$$\Phi(\varepsilon/\sqrt{2L^*}) \geq 1/2,$$

$\varepsilon_0 = \sqrt{2 \ln 2} \approx 1.18$ is a valid choice. \square

It is necessary to have some condition such as (1) that ensures that labels are far from being uniformly chosen.

Proposition 1. Assume that $\mu(A_{nj}) = 2v(A_{nj})$ for all cells A_{nj} of the partition, then

$$\sqrt{n}\mathbf{E}(L_n - \hat{L}_n) \geq \frac{1}{4}\sqrt{m_n}.$$

Thus in that case $\sqrt{n}\mathbf{E}(L_n - \hat{L}_n) \rightarrow \infty$ as $n \rightarrow \infty$.

Proof. Since $L_n \geq R_n^*$ and under the conditions of the proposition $R_n^* = 1/2$ we have to show that

$$\mathbf{E}\hat{L}_n \leq \frac{1}{2} - \frac{1}{4}\sqrt{\frac{m_n}{n}}.$$

$$\begin{aligned}
\mathbf{E}\hat{L}_n &= \mathbf{E} \sum_{j=1}^{m_n} \min\{v_n(A_{nj}), \mu_n(A_{nj}) - v_n(A_{nj})\} \\
&= \mathbf{E} \sum_{j=1}^{m_n} \left(\frac{\mu_n(A_{nj})}{2} + \min\left\{v_n(A_{nj}) - \frac{\mu_n(A_{nj})}{2}, \frac{\mu_n(A_{nj})}{2} - v_n(A_{nj})\right\} \right) \\
&= \mathbf{E} \sum_{j=1}^{m_n} \left(\frac{\mu_n(A_{nj})}{2} - \left| v_n(A_{nj}) - \frac{\mu_n(A_{nj})}{2} \right| \right) \\
&= \frac{1}{2} - \mathbf{E} \sum_{j=1}^{m_n} \left| v_n(A_{nj}) - \frac{\mu_n(A_{nj})}{2} \right| \\
&= \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{n} \mathbf{E} \left| nv_n(A_{nj}) - \frac{n\mu_n(A_{nj})}{2} \right| \\
&= \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{n} \sum_{k=0}^n \mathbf{E} \left(\left| nv_n(A_{nj}) - \frac{n\mu_n(A_{nj})}{2} \right| \mid n\mu_n(A_{nj}) = k \right) \mathbf{P}(n\mu_n(A_{nj}) = k).
\end{aligned} \tag{6}$$

If $n\mu_n(A_{nj}) = k$, then $nv_n(A_{nj})$ is binomial with parameters k and $\frac{1}{2}$, therefore we are interested in the expected value $\mathbf{E}|\text{Binom}(k, \frac{1}{2}) - \frac{k}{2}|$. Let Z_i , $i = 1, \dots, k$ be random variables with distribution $\mathbf{P}(Z_i = 1) = \mathbf{P}(Z_i = -1) = \frac{1}{2}$. Then $|\text{Binom}(k, \frac{1}{2}) - \frac{k}{2}| = \frac{1}{2} |\sum_{i=1}^k Z_i|$ thus applying Lemma A.4 to these Z_i and $p = 1$, we get

$$\frac{1}{2\sqrt{2}} \sqrt{k} \leq \mathbf{E} \left| \text{Binom} \left(k, \frac{1}{2} \right) - \frac{k}{2} \right| = \frac{1}{2} \mathbf{E} \left| \sum_{i=1}^k Z_i \right| \leq \frac{1}{2} \sqrt{k}.$$

Thus from (6)

$$\begin{aligned}
\mathbf{E}\hat{L}_n &\leq \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{n} \sum_{k=0}^n \frac{1}{2\sqrt{2}} \sqrt{k} \binom{n}{k} \mu(A_{nj})^k (1 - \mu(A_{nj}))^{n-k} \\
&= \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{n} n\mu(A_{nj}) \sum_{k=1}^n \frac{1}{2\sqrt{2}} \frac{1}{\sqrt{k}} \binom{n-1}{k-1} \mu(A_{nj})^{k-1} (1 - \mu(A_{nj}))^{(n-1)-(k-1)} \\
&= \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{2\sqrt{2}} \frac{1}{n} n\mu(A_{nj}) \sum_{k=0}^{n-1} \frac{1}{\sqrt{k+1}} \binom{n-1}{k} \mu(A_{nj})^k (1 - \mu(A_{nj}))^{n-1-k} \\
&\leq \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{2\sqrt{2}} \mu(A_{nj}) \frac{1}{\sqrt{\sum_{k=0}^{n-1} (k+1) \binom{n-1}{k} \mu(A_{nj})^k (1 - \mu(A_{nj}))^{n-1-k}}} \\
&\quad \text{(from Jensen's inequality)} \\
&= \frac{1}{2} - \sum_{j=1}^{m_n} \frac{1}{2\sqrt{2}} \mu(A_{nj}) \frac{1}{\sqrt{(n-1)\mu(A_{nj}) + 1}} \\
&\leq \frac{1}{2} - m_n \frac{1}{2\sqrt{2}} \frac{\frac{1}{m_n} \sum_{j=1}^{m_n} \mu(A_{nj})}{\sqrt{(n-1)\frac{1}{m_n} \sum_{j=1}^{m_n} \mu(A_{nj}) + 1}}
\end{aligned}$$

$$\begin{aligned}
& \text{(from Jensen's inequality)} \\
&= \frac{1}{2} - \frac{1}{2\sqrt{2}} \frac{1}{\sqrt{(n-1)/m_n + 1}} \\
&\leq \frac{1}{2} - \frac{1}{4} \sqrt{\frac{m_n}{n}}. \quad \square
\end{aligned}$$

3.2. Deleted estimate

The deleted estimate or cross-validation attempts to avoid the bias present in the resubstitution estimate. The method deletes (X_i, Y_i) from the training data and creates a classifier g_{n-1} using the remaining $n-1$ pairs. It tests for an error on (X_i, Y_i) and repeats this procedure for all n pairs of the training data D_n . Formally denote the training set with (X_i, Y_i) deleted by

$$D_{n,i} = ((X_1, Y_1), \dots, (X_{i-1}, Y_{i-1}), (X_{i+1}, Y_{i+1}), \dots, (X_n, Y_n)).$$

Then define

$$\tilde{L}_n = \frac{1}{n} \sum_{i=1}^n I_{\{g_{n-1}(X_i, D_{n,i}) \neq Y_i\}}.$$

Clearly, the deleted estimate is almost unbiased in the sense that

$$\mathbf{E} \tilde{L}_n = \mathbf{E} L_{n-1}.$$

Let g_{n-1} be the partitioning classification rule for $n-1$ samples, using partition \mathcal{P}_n . Concerning the deleted error estimate for partitioning rule the following is known (see [10, Section 24.5]):

$$\mathbf{E} \{(\tilde{L}_n - L_n)^2\} \leq \frac{1 + 6/e}{n} + \frac{6}{\sqrt{\pi(n-1)}}.$$

The resubstitution estimate for any partitioning classification rule is smaller than the deleted estimate:

$$\tilde{L}_n \geq \hat{L}_n, \quad (7)$$

since if (X_i, Y_i) is a mistake w.r.t. \hat{L}_n , i.e. $Y_i \neq g_n(X_i)$, then the label Y_i is in the minority among the labels of the data falling into its cell. Then, of course, Y_i is in the minority of its cell w.r.t. $D_{n,i} = D_n \setminus \{X_i, Y_i\}$, which implies that (X_i, Y_i) is a mistake w.r.t. \tilde{L}_n .

The deleted estimate for partitioning classification rule is asymptotically normal under similar conditions as the resubstitution estimate:

Theorem 4. *Under the conditions of Theorem 2,*

$$n^{1/2}(L_n - \tilde{L}_n)/\sqrt{L^*/2} \xrightarrow{\mathcal{D}} N(0, 1).$$

Proof. It is easy to see that if U_n and V_n are random variables and $U_n \xrightarrow{\mathcal{D}} N(0, 1)$ and $V_n \rightarrow 0$ in probability, then $U_n + V_n \xrightarrow{\mathcal{D}} N(0, 1)$.

So, since the conditions are the same as for Theorem 2, it suffices to show that

$$\sqrt{n}(\bar{L}_n - \hat{L}_n) \rightarrow 0$$

in probability as $n \rightarrow \infty$.

If for a nonnegative random variable Z_n , $\mathbf{E}Z_n \rightarrow 0$, then $Z_n \rightarrow 0$ in probability. Thus, since from (7) $\bar{L}_n - \hat{L}_n \geq 0$, we have to prove that

$$\sqrt{n}\mathbf{E}(\bar{L}_n - \hat{L}_n) \rightarrow 0$$

as $n \rightarrow \infty$.

(X_i, Y_i) is a mistake w.r.t. \bar{L}_n , i.e. $g_{n-1}(X_i, D_{n,i}) \neq Y_i$ if Y_i is in the minority of its cell w.r.t. $D_{n,i} = D_n \setminus \{X_i, Y_i\}$. If (X_i, Y_i) is a mistake, then each pair with the same label in that cell is also a mistake. The number of ones in a cell is $nv_n(A_{nj})$ and the number of zeros is $n(\mu_n(A_{nj}) - v_n(A_{nj}))$. Thus, for the partitioning classification rule the deleted estimate can be written in the following form:

$$\begin{aligned} \bar{L}_n &= \sum_{j=1}^{m_n} (v_n(A_{nj}) I_{\{nv_n(A_{nj})-1 \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}} \\ &\quad + (\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj}))-1 < nv_n(A_{nj})\}}). \end{aligned}$$

The resubstitution estimate can be written in a similar form using (4)

$$\begin{aligned} \hat{L}_n &= \sum_{j=1}^{m_n} (v_n(A_{nj}) I_{\{nv_n(A_{nj}) \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}} \\ &\quad + (\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj})) < nv_n(A_{nj})\}}). \end{aligned}$$

Therefore

$$\begin{aligned} \sqrt{n}\mathbf{E}(\bar{L}_n - \hat{L}_n) &= \sqrt{n}\mathbf{E} \left\{ \sum_{j=1}^{m_n} (v_n(A_{nj}) I_{\{nv_n(A_{nj})-1 \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}} \right. \\ &\quad + (\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj}))-1 < nv_n(A_{nj})\}} \\ &\quad - (v_n(A_{nj}) I_{\{nv_n(A_{nj}) \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}} \\ &\quad + (\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj})) < nv_n(A_{nj})\}})) \Big\} \\ &= \sqrt{n}\mathbf{E} \left\{ \sum_{j=1}^{m_n} (v_n(A_{nj}) (I_{\{nv_n(A_{nj})-1 \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}} - I_{\{nv_n(A_{nj}) \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}}) \right. \\ &\quad + (\mu_n(A_{nj}) - v_n(A_{nj})) (I_{\{n(\mu_n(A_{nj})-v_n(A_{nj}))-1 < nv_n(A_{nj})\}} \\ &\quad - I_{\{n(\mu_n(A_{nj})-v_n(A_{nj})) < nv_n(A_{nj})\}})) \Big\} \\ &= \sqrt{n}\mathbf{E} \left\{ \sum_{j=1}^{m_n} (v_n(A_{nj}) I_{\{nv_n(A_{nj})-1 \leq n(\mu_n(A_{nj})-v_n(A_{nj})) < nv_n(A_{nj})\}} \right. \\ &\quad + (\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj}))-1 < nv_n(A_{nj}) \leq n(\mu_n(A_{nj})-v_n(A_{nj}))\}}) \Big\} \end{aligned}$$

$$\begin{aligned}
&= \sqrt{n} \mathbf{E} \left\{ \sum_{j=1}^{m_n} (v_n(A_{nj}) I_{\{nv_n(A_{nj})=n(\mu_n(A_{nj})-v_n(A_{nj}))+1\}} \right. \\
&\quad \left. + (\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj}))=nv_n(A_{nj})\}}) \right\} \\
&= \sqrt{n} \sum_{j=1}^{m_n} (\mathbf{E}\{v_n(A_{nj}) I_{\{nv_n(A_{nj})=n(\mu_n(A_{nj})-v_n(A_{nj}))+1\}}\} \\
&\quad + \mathbf{E}\{(\mu_n(A_{nj}) - v_n(A_{nj})) I_{\{n(\mu_n(A_{nj})-v_n(A_{nj}))=nv_n(A_{nj})\}}\}) \\
&= \sqrt{n} \sum_{j=1}^{m_n} (T_{1,nj} + T_{2,nj}). \tag{8}
\end{aligned}$$

Consider the first term of the sum for the cell $A = A_{nj}$:

$$\begin{aligned}
T_{1,n} &= \mathbf{E}\{v_n(A) I_{\{nv_n(A)=n(\mu_n(A)-v_n(A))+1\}}\} \\
&= \sum_{k=0}^n \mathbf{E}\{v_n(A) I_{\{nv_n(A)=n(\mu_n(A)-v_n(A))+1\}} | n\mu_n(A) = k\} \mathbf{P}\{n\mu_n(A) = k\} \\
&= \sum_{l=0}^{n/2} \mathbf{E}\{v_n(A) I_{\{nv_n(A)=n(\mu_n(A)-v_n(A))+1\}} | n\mu_n(A) = 2l + 1\} \mathbf{P}\{n\mu_n(A) = 2l + 1\}.
\end{aligned}$$

Since if the number of pairs in a cell is even, then it cannot contain one more zeros than ones, so the indicator is 0. $n\mu_n(A)$ is binomial with parameters n and $\mu(A)$, so

$$\mathbf{P}\{n\mu_n(A) = 2l + 1\} = \binom{n}{2l + 1} \mu(A)^{2l+1} (1 - \mu(A))^{n-(2l+1)}.$$

If $\mu_n(A) = 2l + 1$, then $nv_n(A)$ is binomial with parameters $2l + 1$ and $v(A)/\mu(A)$, thus

$$\begin{aligned}
&\mathbf{E}\{v_n(A) I_{\{nv_n(A)=n(\mu_n(A)-v_n(A))+1\}} | n\mu_n(A) = 2l + 1\} \\
&= \mathbf{E}\{v_n(A) I_{\{nv_n(A)=l+1\}}\} \\
&= \frac{l+1}{n} \mathbf{P}\{nv_n(A) = l + 1\} \\
&= \frac{l+1}{n} \binom{2l+1}{l+1} \left(\frac{v(A)}{\mu(A)}\right)^{l+1} \left(\frac{\mu(A)-v(A)}{\mu(A)}\right)^l \\
&= \frac{2l+1}{n} \binom{2l}{l} \left(\frac{v(A)}{\mu(A)}\right)^{l+1} \left(\frac{\mu(A)-v(A)}{\mu(A)}\right)^l.
\end{aligned}$$

Therefore

$$\begin{aligned}
T_{1,n} &= \sum_{l=0}^{n/2} \frac{2l+1}{n} \binom{2l}{l} v(A)^{l+1} (\mu(A) - v(A))^l \frac{1}{\mu(A)^{2l+1}} \\
&\quad \times \binom{n}{2l+1} \mu(A)^{2l+1} (1 - \mu(A))^{n-(2l+1)} \\
&= \sum_{l=0}^{n/2} \frac{2l+1}{n} \binom{2l}{l} v(A)^{l+1} (\mu(A) - v(A))^l \binom{n}{2l+1} (1 - \mu(A))^{n-(2l+1)}
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{l=0}^{n/2} 2^{2l} v(A)^{l+1} (\mu(A) - v(A))^l \binom{n}{2l} (1 - \mu(A))^{n-(2l+1)} \\
&\quad \left(\text{since } \frac{2l+1}{n} \binom{2l}{l} \binom{n}{2l+1} = \binom{2l}{l} \binom{n-1}{2l} \leq 2^{2l} \binom{n}{2l} \right) \\
&= \frac{v(A)}{1 - \mu(A)} \sum_{l=0}^{n/2} 2^{2l} v(A)^l (\mu(A) - v(A))^l \binom{n}{2l} (1 - \mu(A))^{n-2l} \\
&= \frac{v(A)}{1 - \mu(A)} \sum_{l=0}^{n/2} \binom{n}{2l} (2\sqrt{v(A)(\mu(A) - v(A))})^{2l} (1 - \mu(A))^{n-2l} \\
&\leq \frac{v(A)}{1 - \mu(A)} \sum_{k=0}^n \binom{n}{k} (2\sqrt{v(A)(\mu(A) - v(A))})^k (1 - \mu(A))^{n-k} \\
&= \frac{v(A)}{1 - \mu(A)} (2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n \\
&\leq \frac{\mu(A)}{1 - \mu(A)} (2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n. \tag{9}
\end{aligned}$$

Similarly, we can get the same upper bound for the second term in (8):

$$\begin{aligned}
T_{2,n} &= \mathbf{E}\{(\mu_n(A) - v_n(A)) I_{\{\mu_n(A) - v_n(A) = nv_n(A)\}}\} \\
&= \sum_{k=0}^n \mathbf{E}\{(\mu_n(A) - v_n(A)) I_{\{\mu_n(A) - v_n(A) = nv_n(A)\}} | n\mu_n(A) = k\} \mathbf{P}\{n\mu_n(A) = k\} \\
&= \sum_{l=0}^{n/2} \mathbf{E}\{(\mu_n(A) - v_n(A)) I_{\{\mu_n(A) - v_n(A) = nv_n(A)\}} | n\mu_n(A) = 2l\} \mathbf{P}\{n\mu_n(A) = 2l\} \\
&= \sum_{l=0}^{n/2} \frac{l}{n} \binom{2l}{l} \left(\frac{v(A)}{\mu(A)}\right)^l \left(\frac{\mu(A) - v(A)}{\mu(A)}\right)^l \binom{n}{2l} \mu(A)^{2l} (1 - \mu(A))^{n-2l} \\
&= \sum_{l=0}^{n/2} \frac{l}{n} \binom{2l}{l} v(A)^l (\mu(A) - v(A))^l \binom{n}{2l} (1 - \mu(A))^{n-2l} \\
&\leq \frac{1}{n} \sum_{l=0}^{n/2} 2l 2^{2l} v(A)^l (\mu(A) - v(A))^l \binom{n}{2l} (1 - \mu(A))^{n-2l} \\
&= \frac{1}{n} \sum_{l=0}^{n/2} 2l \binom{n}{2l} (2\sqrt{v(A)(\mu(A) - v(A))})^{2l} (1 - \mu(A))^{n-2l} \\
&\leq \frac{1}{n} \sum_{k=0}^n k \binom{n}{k} (2\sqrt{v(A)(\mu(A) - v(A))})^k (1 - \mu(A))^{n-k}.
\end{aligned}$$

The sum in the right-hand side is the expected value of the binomial distribution with parameters n and $2\sqrt{v(A)(\mu(A) - v(A))}/(2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))$ if we normalize with the sum of $2\sqrt{v(A)(\mu(A) - v(A))}$ and $1 - \mu(A)$. Thus

$$\begin{aligned}
T_{2,n} &\leq \frac{1}{n} (2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n \frac{1}{(2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n} \\
&\quad \times \sum_{k=0}^n k \binom{n}{k} (2\sqrt{v(A)(\mu(A) - v(A))})^k (1 - \mu(A))^{n-k}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} (2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n \\
&\quad \times \mathbf{E} \left\{ \text{Binom} \left(n, \frac{2\sqrt{v(A)(\mu(A) - v(A))}}{2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A)} \right) \right\} \\
&= \frac{1}{n} \frac{n 2\sqrt{v(A)(\mu(A) - v(A))}}{2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A)} (2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n \\
&\leq \frac{\mu(A)}{1 - \mu(A)} (2\sqrt{v(A)(\mu(A) - v(A))} + 1 - \mu(A))^n.
\end{aligned}$$

From this upper bound and from (8) and (9)

$$\begin{aligned}
&\sqrt{n} \mathbf{E}(\tilde{L}_n - \hat{L}_n) \\
&\leq \sqrt{n} \sum_{j=1}^{m_n} 2 \frac{\mu(A_{nj})}{1 - \mu(A_{nj})} (2\sqrt{v(A_{nj})(\mu(A_{nj}) - v(A_{nj}))} + 1 - \mu(A_{nj}))^n \\
&\leq \sqrt{n} \sum_{j=1}^{m_n} 2\mu(A_{nj}) e^{-n(\mu(A_{nj}) - 2\sqrt{v(A_{nj})(\mu(A_{nj}) - v(A_{nj}))})} \\
&= \sqrt{n} \sum_{j=1}^{m_n} 2\mu(A_{nj}) e^{-n\mu(A_{nj})(1 - 2\sqrt{\frac{v(A_{nj})}{\mu(A_{nj})}(1 - \frac{v(A_{nj})}{\mu(A_{nj})})})},
\end{aligned}$$

where the second inequality comes from the fact that $1 - z \leq e^{-z}$ if $z > 0$, and clearly $\mu(A_{nj}) - 2\sqrt{v(A_{nj})(\mu(A_{nj}) - v(A_{nj}))} > 0$, and we can assume that $\mu(A_{nj}) < \frac{1}{2}$.

The condition (1) means that

$$\left| \frac{1}{2} - \frac{v(A_{nj})}{\mu(A_{nj})} \right| > \frac{c}{2},$$

and because of this

$$\frac{v(A_{nj})}{\mu(A_{nj})} \left(1 - \frac{v(A_{nj})}{\mu(A_{nj})} \right) \leq \frac{1}{4} - \frac{c^2}{4}.$$

Thus, if we denote $1 - 2\sqrt{\frac{1}{4} - \frac{c^2}{4}} = 1 - \sqrt{1 - c^2}$ by $\delta(c)$ then

$$\begin{aligned}
\sqrt{n} \mathbf{E}(\tilde{L}_n - \hat{L}_n) &\leq \sqrt{n} \sum_{j=1}^{m_n} 2\mu(A_{nj}) e^{-n\mu(A_{nj})(1 - 2\sqrt{\frac{v(A_{nj})}{\mu(A_{nj})}(1 - \frac{v(A_{nj})}{\mu(A_{nj})})})} \\
&\leq \sqrt{n} \sum_{j=1}^{m_n} 2\mu(A_{nj}) e^{-\delta(c)n\mu(A_{nj})} \\
&= \sqrt{n} \sum_{j=1}^{m_n} \frac{2}{\delta(c)n} \delta(c)n\mu(A_{nj}) e^{-\delta(c)n\mu(A_{nj})} \\
&\leq \sqrt{n} \sum_{j=1}^{m_n} \frac{2}{\delta(c)n} \max_z z e^{-z} \\
&\leq \sqrt{nm_n} \frac{2}{\delta(c)n} = \frac{2m_n}{\delta(c)\sqrt{n}} \rightarrow 0
\end{aligned}$$

because of condition (3). \square

Thus, similar to the resubstitution estimate, the deleted estimate is a good estimate of the true error probability of the partitioning classification rule, and using the result of the theorem a good estimate can be given to the probability that the difference between the deleted estimate and the error probability is greater than ε/\sqrt{n} .

Appendix

Lemma A.1 (Bernstein [6]). *For independent identically distributed zero mean random variables Z_i with $|Z_i| \leq t$, and for all $\varepsilon > 0$*

$$\mathbf{P}\left(\left|\frac{1}{n} \sum_{i=1}^n Z_i\right| > \varepsilon\right) \leq 2 \exp\left(\frac{-n\varepsilon^2}{2(\sigma^2 + t\varepsilon)}\right),$$

where $\sigma^2 = \mathbf{E}(Z_i^2)$.

Lemma A.2 (Györfi and Horváth [14]). *Under the conditions of Theorem 2,*

$$n^{1/2} (\hat{L}_n - R_n^*) / \sqrt{L^*/2} \xrightarrow{\mathcal{D}} N(0, 1).$$

Lemma A.3 (McDiarmid [15]). *Let Z_1, \dots, Z_n be independent random variables taking values in a set A , and assume that the measurable function $F: A^n \rightarrow \mathcal{R}$ satisfies*

$$\sup_{z_1, \dots, z_n, z'_i} |F(z_1, \dots, z_i, \dots, z_n) - F(z_1, \dots, z'_i, \dots, z_n)| \leq c_i, \quad 1 \leq i \leq n.$$

Then for all $\varepsilon > 0$

$$\mathbf{P}\{|F(Z_1, \dots, Z_n) - \mathbf{E}F(Z_1, \dots, Z_n)| > \varepsilon\} \leq 2e^{-2\varepsilon^2/(\sum_{i=1}^n c_i^2)}.$$

Lemma A.4 (Khinchine's inequality, see e.g. [9]). *Let Z_1, \dots, Z_k be independent random variables taking the values $+1$ and -1 with equal probability, and let a_1, \dots, a_k be real numbers. Then there exist positive constants B_p and C_p depending upon $p > 0$ only such that*

$$B_p \sqrt{\sum_{i=1}^k a_i^2} \leq \mathbf{E}^{1/p} \left(\left| \sum_{i=1}^k a_i Z_i \right|^p \right) \leq C_p \sqrt{\sum_{i=1}^k a_i^2}.$$

The following values are optimal:

$$B_p = \begin{cases} 2^{1/2-1/p} & \text{if } 0 < p \leq p_0, \\ 2^{1/2} \left(\frac{\Gamma((p+1)/2)}{\sqrt{\pi}} \right)^{1/p}, & p_0 \leq p < 2, \\ 1 & \text{if } 2 \leq p, \end{cases}$$

$$C_p = \begin{cases} 1 & \text{if } 0 < p \leq 2, \\ 2^{1/2} \left(\frac{\Gamma((p+1)/2)}{\sqrt{\pi}} \right)^{1/p}, & 2 \leq p, \end{cases}$$

where $\Gamma(\cdot)$ is the classical gamma-function extending the factorial and $p_0 = 1.84742 \dots$ is the solution in $(1, 2)$ of $\Gamma((p+1)/2) = \Gamma(\frac{3}{2})$.

References

- [1] M.S. Bartlett, The characteristic function of a conditional statistic, *J. London Math. Soc.* 13 (1938) 62–67.
- [2] J. Beirlant, L. Györfi, On the asymptotic normality of L_2 -error in partitioning regression estimation, *J. Statist. Plann. Inference* 71 (1998) 93–107.
- [3] J. Beirlant, L. Györfi, G. Lugosi, On the asymptotic normality of L_1 - and L_2 -errors in histogram density estimation, *Canad. J. Statist.* 22 (3) (1994) 309–318.
- [4] J. Beirlant, D.M. Mason, On the asymptotic normality of L_p -norms of empirical functionals, *Math. Methods Statist.* 4 (1) (1995) 1–19.
- [5] A. Berline, L. Devroye, L. Györfi, Asymptotic normality of L_1 -error in histogram density estimation, *Statistics* 26 (1995) 329–343.
- [6] S. Bernstein, *The Theory of Probabilities*, Gostehizdat Publishing House, Moscow, 1946.
- [7] L. Devroye, Necessary and sufficient conditions for the almost everywhere convergence of the nearest neighbor regression function estimates, *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete* 61 (1982) 467–481.
- [8] L. Devroye, L. Györfi, Distribution-free exponential bound on the L_1 error of partitioning estimates of a regression function, in: F. Konecny, J. Mogyoródi, W. Wertz (Eds.), *Proceedings of the Fourth Pannonian Symposium on Mathematical Statistics*, Akadémiai Kiadó, Budapest, Hungary, 1983, pp. 67–76.
- [9] L. Devroye, L. Györfi, *Nonparametric Density Estimation: The L_1 View*, John Wiley, New York, 1985.
- [10] L. Devroye, L. Györfi, G. Lugosi, *Probabilistic Theory of Pattern Recognition*, Springer, New York, 1996.
- [11] N. Glick, Sample based multinomial classification, *Biometrics* 29 (1973) 241–256.
- [12] L. Gordon, R. Olshen, Asymptotically efficient solutions to the classification problem, *Ann. Statist.* 6 (1978) 515–533.
- [13] L. Gordon, R. Olshen, Consistent nonparametric regression from recursive partitioning schemes, *J. Multivariate Anal.* 10 (1980) 611–627.
- [14] L. Györfi, M. Horváth, On the asymptotic normality of the resubstitution error estimate for partitioning classification rule, in: A. Rizza, M. Vichi, H.H. Bock (Eds.), *Advances in Data Science and Classification*, Springer, Berlin, 1998, pp. 197–204.
- [15] C. McDiarmid, On the method of bounded differences, *Surveys in Combinatorics 1989*, Cambridge University Press, Cambridge, 1989, pp. 148–188.
- [16] J.W. Tukey, Nonparametric estimation. II. Statistically Equivalent Blocks and Tolerance Regions, *Ann. Math. Statist.* 18 (1947) 529–539.